

Evaluating large language models for accuracy incentivizes hallucinations

<https://doi.org/10.1038/s41586-026-10549-w>

Adam Tauman Kalai¹, Ofir Nachum¹, Santosh S. Vempala² & Edwin Zhang^{1,3}

Received: 1 July 2025

Accepted: 15 April 2026

Published online: 22 April 2026

Open access

 Check for updates

Large language models sometimes produce confident, plausible falsehoods ('hallucinations'), limiting their reliability^{1,2}. Previous work has offered numerous explanations and effective mitigations such as retrieval and tool use³, consistency-based self-verification⁴ and reinforcement learning from human feedback⁵. Nonetheless, the problem persists even in state-of-the-art language models^{6,7}. Here we show how next-word prediction and accuracy-based evaluations inadvertently reward unwarranted guessing. Initially, next-word pretraining creates statistical pressure towards hallucination even with idealized error-free data: using learning theory^{8,9}, we show that facts lacking repeated support in training data (such as one-off details) yield unavoidable errors, whereas recurring regularities (such as grammar) do not. Subsequent training stages aim to correct such errors. However, dominant headline metrics such as accuracy systematically reward guessing over admitting uncertainty. To align incentives, we suggest two additions to the classic approach of adding error penalties to evaluations to control abstention^{10,11}. First, we propose 'open rubric' evaluations that explicitly state how errors are penalized (if at all), which test whether a model modulates its abstentions to stated stakes while optimizing accuracy. Second, as hallucination-specific benchmarks rarely make leaderboards¹², we suggest using open-rubric variants of existing evaluations, to reverse their guessing incentives. Reframing hallucination as an incentive problem opens a practical path towards more reliable language models.

Plausible, confidently stated falsehoods diminish the utility of large language models (LLMs) in reliability-critical domains. Despite progress, this problem persists even in state-of-the-art models^{6,7}. For example, three popular LLMs (details in Methods) responded to the question 'What does PGGB stand for?' as follows:

ChatGPT: 'PGGB' can stand for different things depending on context, but one of the most common meanings is 'Polynomial Gaussian Gradient Bandwidth'...

Claude: PGGB most commonly stands for Privately Held Global Growth and Income Fund, a closed-end fund that trades on the NYSE...

DeepSeek: Of course! PGGB stands for Perfect Guard Group Buy...

All are wrong. They do not abstain, for example, by saying, 'I don't know', or by requesting more context, as a reliable human assistant would. Instead, they fabricate specific, confident responses. Like students facing a difficult exam question, models guess rather than admit uncertainty. To reduce hallucinations further, we need to understand why they arise and persist (Fig. 1).

Various explanations have been offered for how hallucinations arise¹; here we provide a unifying lens using computational learning

theory⁸. To make this lens precise, we first address a practical obstacle: defining and measuring hallucination is complicated by issues such as responses that contain multiple vague claims¹³. We side-step these definitional tangles by analysing an abstract set of errors, following learning theory, which applies to binary classification (for example, dog versus cat images) without adjudicating edge cases (for example, images containing both).

Initially, LLMs are trained to optimize next-word (or next-token) prediction, during pretraining. Although falsehoods in training data can cause hallucinations, the phenomenon is not purely garbage-in garbage-out. We show that this objective creates statistical tendency towards hallucination even with ideal error-free training data. The key insight is a formal reduction to binary classification: any language model implicitly answers 'Is this response valid?' for each candidate generation. This connection establishes lower bounds on hallucination rates and illuminates which error types to expect, as causes of misclassification are well understood⁹. Learnable patterns (grammar, spelling, politeness) yield low error rates, whereas non-learnable facts (birthdays, one-off details) produce high hallucination rates in pre-trained LLMs. A natural conjecture is that removing errors and including abstentions in the training corpus would resolve most hallucinations, but our analysis shows that this is not the case.

Hallucinations persist, despite the fact that numerous alignment approaches have been shown to curb them²⁻⁵. We offer a simple explanation: under standard scoring, guessing is a dominant strategy. In terms

¹OpenAI, San Francisco, CA, USA. ²Georgia Institute of Technology, Atlanta, GA, USA. ³Present address: Isara Laboratories, San Francisco, CA, USA. ✉e-mail: adam@ka.lai

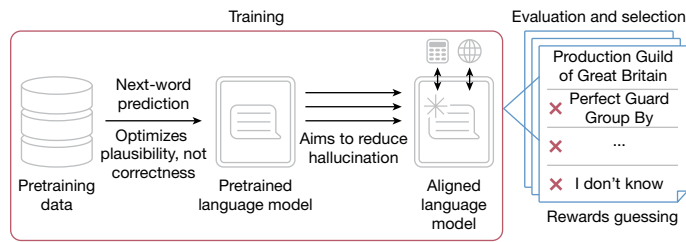


Fig. 1 | Origins and persistence of hallucination across training and evaluation. Next-word prediction creates statistical pressure towards hallucination. Subsequently, numerous techniques improve alignment and reduce hallucination. Yet accuracy-based evaluations ultimately inform model selection, inadvertently incentivizing confident guessing.

of raw accuracy or pass rate (fraction of overall questions answered correctly), a model that guesses outscore a trustworthy model that abstains when uncertain. For example, on the SimpleQA evaluation¹⁴, accuracy slightly favours OpenAI’s o4-mini, which answers almost all questions (with over 3/4 error rate)¹⁵ over GPT-5-mini, even though GPT-5-mini makes many fewer errors (owing to abstentions)⁶. Our meta-evaluation of popular LLM benchmarks confirms that the majority favour guessing. The ubiquity of rewarding guessing explains why the handful of existing hallucination evaluations also have not resolved the problem¹². Instead, we argue that existing primary evaluations need to be revised. One may hope that accuracy will approach 100% through innovations such as web search and scaling. However, some real-world questions will remain unanswerable (for example, unlisted birthdays or impossible questions), so accuracy-style metrics still favour guessing when uncertain.

Finally, evaluation incentives need to be realigned to reduce hallucination. A key challenge is that the true response utility is context specific and subjective. OpenAI¹⁶ specifies only a vague ranking: correct answer > no answer > wrong answer. If numeric scores are not aligned with utility, overconfident or underconfident models will outscore the ideal LLM that appropriately expresses uncertainty.

We propose a straightforward method to modify existing questions to make scoring unambiguous: state the scoring system alongside the question itself. Such ‘open rubric’ questions can be paired with the classic approach of encouraging abstention by penalizing errors^{10,11}. For example, for the recent ‘net score’¹⁷, prompts could be augmented with ‘correct answers receive 1 and incorrect answers receive -1 (hence abstain if <50% likely to be correct)’. When there is no penalty, ‘only fully correct answers receive credit (so make your best guess if unsure)’ describes standard accuracy. Under open rubrics, accuracy is no longer at odds with humility: a reliable model can guess when instructed to do so and abstain when not, just as people guess more on exams than in everyday settings. The GPT-5-mini model mentioned above, which often abstained, now becomes more accurate than o4-mini when measured with open rubrics (Methods). A model that performs well across penalties demonstrates controllable abstention. We provide a case study using four frontier models, demonstrating how open rubrics incentivize adoption of a hallucination-reduction technique.

Taken together, our analysis reframes hallucination as an unintended outcome of training objectives and evaluation incentives rather than an inherent LLM deficiency, and demonstrates how aligned evaluations can incentivize reliability in language models.

Hallucination from next-word pretraining

Our analysis offers an organizing lens for a range of error phenomena, including hallucination, that arise when generating with calibrated next-word predictors. This perspective clarifies how hallucinations differ statistically from more regular errors such as misspellings.

Pretraining fits a probability model \hat{p} to text sampled from a distribution p using maximum-likelihood estimation, as in autocomplete systems. With ideal training data, autocomplete should generate correct spelling, yet it may still fabricate dates when completing ‘Adam Tauman Kalai was born on...’ if that birthday is absent from the training data. (This is only expected from pretraining, before further alignment steps.) The challenge in formalizing this intuition is that falsehoods would not be generated by the perfect predictor $\hat{p} = p$ trained on infinite error-free training data; nor would they be generated by the untrained LLM that completes everything with ‘I don’t know’. Our analysis applies to large but finite error-free datasets that may contain ample abstentions. With errors in training data, one may expect a higher hallucination rate.

We show that generating valid outputs, for calibrated LLMs, is harder than classifying output validity. Specifically, consider an is-it-valid (IIV) binary-classification problem that has a training set consisting of a large number of responses, each labelled either as valid (+) or error (-), as illustrated in Fig. 2. For this supervised learning problem, both train and test data are 50/50 mixtures of valid examples labelled as + (that is, the pretraining data as we assume it is valid) and uniformly random errors labelled as -. We then show how any LLM can be used as an IIV classifier. Theorem 1 relates IIV and generative errors (such as hallucinations):

$$(\text{Generative error rate}) \geq 2 \times (\text{IIV misclassification rate}).$$

Well-understood error decompositions in supervised learning⁹ thus translate to generative errors. Figure 2 (right) illustrates two causes visually: middle, a poor model of a linear separator for a circular region (that is, approximation error owing to hypothesis-class misspecification); and bottom, statistical complexity (that is, estimation error owing to finite samples and capacity control).

Statistical complexity is a common cause of hallucination arising when there is no succinct pattern that fully describes the concept. Knowledge gaps are inevitable with non-exhaustive training data, for example, one cannot predict birthdays absent from the training data. Using the IIV reduction, Theorem 3 (Methods) shows a stylized setting where

$$(\text{Hallucination rate}) \geq (\text{Fraction of training facts that appear exactly once}).$$

For instance, if 20% of birthday facts appear exactly once in pretraining data, then pretrained models should hallucinate on at least 20% of birthday facts. However, LLMs rarely hallucinate on country capitals, which all appear repeatedly in training data. Theorem 3 extends previous work in this stylized setting¹⁷ that did not consider prompts or abstention.

We now present our main reduction without prompts. The analysis with prompts (Theorem 2 in Methods) is similar. Without prompts, a pretrained LLM \hat{p} is a probability distribution over a set \mathcal{X} of ‘plausible’ texts x , for example, statements or documents. We assume that \mathcal{X} is finite for simplicity. The examples $\mathcal{X} = \mathcal{E} \cup \mathcal{V}$ are partitioned into errors \mathcal{E} and valid examples \mathcal{V} , for non-empty disjoint sets \mathcal{E}, \mathcal{V} . For hallucinations, \mathcal{E} would consist of plausible generations containing one or more falsehoods. The error rate of LLM \hat{p} is

$$\text{err} := \hat{p}(\mathcal{E}) = \Pr_{x \sim \hat{p}} [x \in \mathcal{E}].$$

Here and throughout, $x \sim q$ means that x is distributed according to q . Pretraining data are assumed to come from a noiseless pretraining distribution p over \mathcal{X} , that is, $p(\mathcal{E}) = 0$.

The IIV binary-classification problem is formally specified by the target function $f: \mathcal{X} \rightarrow \{-, +\}$ to be learned (membership in \mathcal{V}) and the distribution D over examples \mathcal{X} (a 50/50 mix of samples from p and uniformly random errors):

Valid examples +		Error examples -	
Greetings. How can I help?	Greetings. How kan eye help?		Spelling (good model)
There are 2 Ds in LADDER There is 1 N in PIANO	There are 3 Ls in SPELL There is 1 G in CAT		Counting (poor model)
Mia Holdner's birthday is 4/1. I don't know Zdan's birthday.	Colin Merivale's birthday is 8/29. Jago Pere's birthday is 8/21.		Birthdays (no pattern)

Fig. 2 | Generative errors as IIV misclassification. Left: IIV requires learning to identify valid generations using labelled valid (+) and error (-) examples. Right: each row shows a concept and a learned classifier (dashed line); spelling is separable with a good model (top), counting errors are due to a poor model representation (middle), and arbitrary facts have errors because there is no pattern in the data (bottom).

$$D(x) := \begin{cases} p(x)/2 & \text{if } x \in \mathcal{V}, \\ 1/(2|\mathcal{E}|) & \text{if } x \in \mathcal{E}, \end{cases} \text{ and } f(x) := \begin{cases} + & \text{if } x \in \mathcal{V}, \\ - & \text{if } x \in \mathcal{E}. \end{cases}$$

The LLM is thus used as an IIV classifier, in our reduction, by thresholding its probability $\hat{p}(x)$ (which can generally be efficiently computed) at $1/|\mathcal{E}|$. The IIV misclassification rate is

$$\text{err}_{\text{iiv}} := \Pr_{x \sim D} [\hat{f}(x) \neq f(x)], \text{ where } \hat{f}(x) := \begin{cases} + & \text{if } \hat{p}(x) > 1/|\mathcal{E}|, \\ - & \text{if } \hat{p}(x) \leq 1/|\mathcal{E}|. \end{cases}$$

Theorem 1. For any pretraining distribution p such that $p(\mathcal{E}) = 0$ and any LLM \hat{p}

$$\text{err} \geq 2 \times \text{err}_{\text{iiv}} - \frac{|\mathcal{V}|}{|\mathcal{E}|} - \delta,$$

where $\delta := |\hat{p}(\mathcal{T}) - p(\mathcal{T})|$ and $\mathcal{T} := \{x \in \mathcal{X} \mid \hat{p}(x) > 1/|\mathcal{E}|\}$.

As Theorem 1 holds for any LLM \hat{p} , it immediately implies that all pretrained LLMs will err on inherently unlearnable IIV facts (such as birthdays absent from the training data) where err_{iiv} is necessarily large, and where δ and $|\mathcal{V}|/|\mathcal{E}|$ are small (for example, for each person there are 364 more incorrect birthday claims in \mathcal{E} than correct ones in \mathcal{V} , although \mathcal{V} also includes valid abstentions).

As pretraining aims for $\hat{p} \approx p$, one may expect small $\delta = |\hat{p}(\mathcal{T}) - p(\mathcal{T})|$, which is a weak form of what is known as calibration. Formally, it has been shown that maximum-likelihood estimation (that is, minimizing cross-entropy) leads to calibration¹⁸, and pretrained LLMs empirically satisfy related forms of calibration^{19,20}. However, aligned models should not hallucinate and hence should be poorly calibrated. Methods contains proofs and shows how other types of hallucination are related to misclassification.

Evaluation metrics that reward guessing

Beyond next-word prediction, a number of techniques have demonstrated significant empirical reductions in hallucination rates², some of which show substantial gains even on frontier models. Yet these methods are typically judged (and models are selected) using headline metrics such as accuracy or pass rate, which treat abstention as failure and thus reward guessing. Here we observe that the metrics used in existing benchmarks and leaderboards reinforce hallucination when a model is unsure.

Most leaderboards score each problem (or each subpart) as either correct or incorrect, and abstention is typically graded as incorrect. Such binary evaluations impose a false right-wrong dichotomy, awarding no credit to answers that appropriately express uncertainty, omit dubious details or request clarification. Under this scoring, abstaining is strictly suboptimal, being penalized as incorrect while an overconfident

'best guess' is optimal for maximizing expected accuracy. As a result, accuracy-based evaluation pushes models to convert uncertainty into content that is often plausible but wrong. Observation 1 (Methods) formally encapsulates how guessing is a dominant strategy for binary evaluations. Reducing hallucination often means lower accuracy, which impedes adoption.

Extended Data Table 1 summarizes our meta-evaluation (Methods), finding that the vast majority of popular evaluations use binary grading. Therefore, additional hallucination evaluations may not suffice when the primary evaluations penalize honestly reporting uncertainty. This does not diminish existing work on hallucination evaluations but rather highlights that even an ideal one may still be outweighed by lower scores on the vast majority of existing benchmarks.

Scoring metrics do not merely measure model performance: in practice, scores are effectively optimized throughout LLM development, including the selection of data, architectures and algorithms. As a result, accuracy-centric leaderboards can perpetuate hallucination by rewarding guessing when a model is unsure, motivating evaluation designs that break this cycle.

Aligning evaluation incentives

Because evaluations inform which LLMs are deployed, their scoring shapes model behaviour. A classic approach to discouraging guessing is to penalize errors^{10,11}. More generally, for each response r to a question prompt (or context) $c \in \mathcal{C}$, we seek a score $s(c, r) \in \mathbb{R}$ that approximates response utility $u(c, r)$. Utility is context-specific: confabulated acronyms are less costly than hallucinations about elevator capacity limits, and overestimates on a capacity sign are worse than small mistakes helping students study. The recent 'net score'⁷ assigns 1 to correct answers and -1 to incorrect answers. This incentivizes guessing only when the probability of being correct exceeds 50%, which is not uniformly appropriate (although arguably superior to the status quo 0% threshold).

One way to make scoring unambiguous is to explicitly specify the scoring rubric in the prompt, removing the need to subjectively assess error severity. We refer to a question as open rubric if it is clear to the test-taker how errors and abstentions are scored. Sample instruction templates used in our case study are in Extended Data Table 2.

Error penalties have traditionally been applied at the scoring stage and are therefore invisible to the model. A model that is told the cost of errors can adjust abstention to the stakes, much as students adjust guessing on exams to a known grading system. An error penalty $L \geq 0$ is strategically equivalent to offering a reward of $t := L/(1+L)$ for abstaining. Although error penalties are common in human exams, we use abstention rewards here because incentives are simpler to understand: if abstaining earns 75% credit, then it is optimal to guess when one is more than 75% likely to be correct. Extended Data Table 2 gives two ways to describe rubrics in text (for example, explaining that errors cost 3 versus abstaining receives 75% credit). In either case, we refer to t as the correctness threshold because both views incentivize guessing when one is more than t likely to be correct.

Formally, the abstention reward function assigns scores 1 (correct), t (abstain) and 0 (error):

$$s_t(c, r) := g_c(r) + t \times \mathbb{1}[r \in \mathcal{A}_c]$$

where $g_c(r) \in \{0, 1\}$ grades correctness for question c and response r , \mathcal{A}_c is a set of abstention responses, and the indicator $\mathbb{1}[\phi]$ denotes 1 if predicate ϕ holds and 0 otherwise. The corresponding error-penalty score is a rescaling $s'_t(c, r) := (s_t(c, r) - t)/(1 - t)$, with 1 (correct), 0 (abstain) and $-L$ (error) where $L = t/(1 - t)$. Let c_t denote the question augmented with the scoring rubric. For alignment, assuming utility $u(c_t, r) = s_t(c_t, r)$ is more reasonable than $u(c, r) = s_t(c, r)$ for t unspecified in c . At least this rewards a desirable behaviour: being correct as often

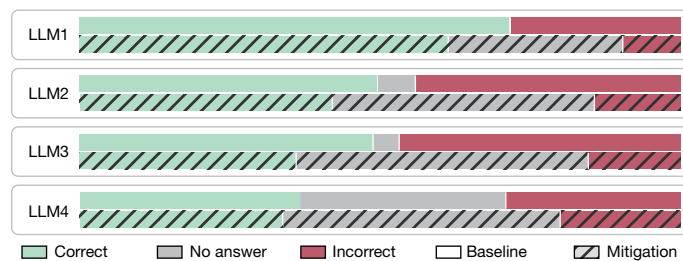


Fig. 3 | Accuracy as a barrier to hallucination reduction. The consistency-based mitigation reduces incorrect responses (hallucination rate) but also decreases correct responses (accuracy), posing a barrier to adoption ($n = 4,326$ questions per model).

as possible while abstaining when the correctness probability is below the stated threshold. Open rubrics mirror typical human settings: students know whether or not they are taking an exam (and how wrong answers are scored), and guess accordingly. By making scoring visible, open rubrics let LLMs make the same distinction rather than collapsing all contexts into a single mode.

To make these incentives concrete, we run a stylized experiment to see whether open-rubric or closed-rubric evaluations are more likely to favour adopting a hallucination mitigation. The mitigation is tested on frontier reasoning models: Google’s Gemini 3 Pro (LLM1), OpenAI’s GPT-5 (LLM2), xAI’s Grok 4 (LLM3) and Anthropic’s Claude Opus 4.5 (LLM4). This is not a controlled evaluation across models: we use default settings, with no tuning or cost normalization.

We use the SimpleQA evaluation¹⁴, which comprises 4,326 factual questions, such as ‘What years was Antonio de Padua María Severino López de Santa Anna y Pérez de Lebrón vice president?’. Accuracy is often used as a headline metric for SimpleQA²¹ despite its authors’ intentions¹⁴. We begin with the standard closed-rubric setting, where the scoring rubric is not mentioned in its questions. The baseline condition prompts the LLM to answer the factual question as stated, without additional instruction. We consider a prototypical hallucination mitigation: (1) the model generates two independent responses; (2) those two responses (in addition to the query) are presented to the same model to judge consistency; and (3) if they are judged to be consistent, the first output is selected, otherwise the model abstains. Similar consistency approaches have been studied for detecting and reducing hallucination^{20,22}. The mitigation was chosen for simplicity; many hallucination-reduction techniques have been studied, some more effective or efficient². As shown in Fig. 3, the mitigation cuts errors but also reduces accuracy across models.

We next append the rubric strings rub, from Extended Data Table 2 to each question with $t = L/(L + 1)$, explicitly stating an abstention reward. (The abstention reward version was chosen for simplicity, but a small ablation with error penalties suggested similar performance.) We use thresholds $t = 0, 0.5, 0.75, 0.9$ corresponding to penalties $L = 0, 1, 3, 9$. Models abstain more for larger thresholds t (Extended Data Fig. 2). We again evaluate baseline and mitigation, now prompting with the augmented question. The same mitigation as above is used (including the appended rubric) for $t > 0$. At $t = 0$, however, because the model is instructed to never abstain, a different prompt is used to select the better of the two responses. Conditioning on different error penalties is impossible with closed rubrics because they are not stated.

With an open rubric, the mitigation outperforms the baseline across models and penalties, including zero (accuracy) (Table 1 and Extended Data Table 3). Open rubrics thus offer consistent encouragement to adopt the mitigation, while closed-rubric accuracy discourages adoption. Open rubrics also provide a natural way to detect overconfidence. For penalties $L = 3, 9$, all four models in Extended Data Table 3 show strong overconfidence with negative (unmitigated) scores, so always abstaining yields a higher score. For open rubrics with no penalty $L = 0$,

Table 1 | Closed rubric@L and open rubric@L mean scores with penalties $L = 0, 1$

	Closed rubric@0	Closed rubric@1	Open rubric@0	Open rubric@1
LLM1	0.71	0.43	0.71	0.46
LLM1 ^a	0.61	0.52	0.73	0.55
LLM2	0.49	0.05	0.50	0.18
LLM2 ^a	0.42	0.28	0.51	0.30
LLM3	0.49	0.02	0.49	0.22
LLM3 ^a	0.36	0.21	0.51	0.26
LLM4	0.37	0.07	0.46	0.14
LLM4 ^a	0.34	0.14	0.47	0.17

When the scoring is unspecified (closed rubric), the mitigation reduces closed rubric@0 (accuracy) for all four LLMs, challenging adoption. With open rubrics, the mitigation helps across the board (for $L = 3, 9$; Extended Data Table 3). Each cell is an average over $n = 4,326$ questions. Values in bold represent a statistically significant difference between the mitigation and baseline ($P < 10^{-5}$ using a two-sided paired permutation test). ^aWith mitigation.

one cannot be overconfident because the instructions are to always make one’s best guess.

Discussion

Current LLM development inadvertently limits reliability through multiple stages: pretraining creates statistical pressure towards hallucination, and the prevalence of closed-rubric accuracy-based evaluations means that LLMs are always in exam mode. Improving reliability is thus not only a modelling problem but also an evaluation mechanism-design problem. Rather than adding separate hallucination benchmarks, we argue that primary evaluations should be modified so that they incentivize admitting uncertainty when appropriate. Open rubrics are one such modification: reducing hallucinations no longer comes at the cost of headline metrics. Error penalties and abstention-aware scoring are established ideas¹⁰, but with LLMs the scoring stakes can be stated explicitly. Evaluating across a range of penalties measures a model’s ability to modulate abstention according to the stated penalty—a good student is both knowledgeable and attuned to when guessing is appropriate.

However, closed rubrics have the advantage of ecological validity, as users rarely specify a confidence threshold in real-world prompts. Furthermore, closed-rubric evaluations with appropriate context-specific scores could measure the ability to assess error severity (although assessing different harms is challenging). The correct–incorrect–abstain trichotomy used here is deliberately simple, and it would be desirable to adapt existing evaluations to handle linguistic calibration (‘I would guess...’)²³ and the broader pragmatics of language^{24,25}. Although hallucination rates have decreased substantially since early models^{6,7}, further progress may depend as much on what we measure as on what we build: aligning evaluation incentives can make reliability improvements pay off.

Online content

Any methods, additional references, Nature Portfolio reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41586-026-10549-w>.

- Huang, L. et al. A survey on hallucination in large language models: principles, taxonomy, challenges, and open questions. *ACM Trans. Inf. Syst.* **43**, 42 (2025).
- Ji, Z. et al. Survey of hallucination in natural language generation. *ACM Comput. Surv.* **55**, 248:1–248:38 (2023).

3. Lewis, P. et al. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *Advances in Neural Information Processing Systems* Vol. 33 (eds Larochelle, H. et al.) 9459–9474 (Curran Associates, 2020); https://proceedings.neurips.cc/paper_files/paper/2020/file/6b493230205f780e1bc26945df7481e5-Paper.pdf.
4. Manakul, P., Liusie, A. & Gales, M. SelfCheckGPT: zero-resource black-box hallucination detection for generative large language models. In *Proc. 2023 Conference on Empirical Methods in Natural Language Processing* (eds Bouamor, H. et al.) 9004–9017 (Association for Computational Linguistics, 2023); <https://aclanthology.org/2023.emnlp-main.557/>.
5. Ouyang, L. et al. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems* Vol. 35 (eds Koyejo, S. et al.) 27730–27744 (Curran Associates, 2022); https://papers.nips.cc/paper_files/paper/2022/hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html.
6. GPT-5 System Card Technical Report (OpenAI, 2025); <https://cdn.openai.com/gpt-5-system-card.pdf>.
7. Claude Opus 4.6 System Card Technical Report (Anthropic PBC, 2026); <https://www-cdn.anthropic.com/Odd865075ad3132672ee0ab40b05a53f14cf5288.pdf>.
8. Kearns, M. J., Schapire, R. E. & Sellie, L. M. Toward efficient agnostic learning. *Mach. Learn.* **17**, 115–141 (1994).
9. Shalev-Shwartz, S. & Ben-David, S. *Understanding Machine Learning: from Theory to Algorithms* (Cambridge Univ. Press, 2014).
10. Chow, C. K. On optimum recognition error and reject tradeoff. *IEEE Trans. Inf. Theory* **16**, 41–46 (1970).
11. Xin, J., Tang, R., Yu, Y. & Lin, J. The art of abstention: selective prediction and error regularization for natural language processing. In *Proc. 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)* (eds Zong, C. et al.) 1040–1051 (Association for Computational Linguistics, 2021); <https://aclanthology.org/2021.acl-long.84/>.
12. Maslej, N. et al. *Artificial Intelligence Index Report 2025* Annual Report (AI Index Steering Committee, Institute for Human-Centered AI, Stanford Univ. 2025); <https://hai.stanford.edu/ai-index/2025-ai-index-report>.
13. van Deemter, K. The pitfalls of defining hallucination. *Comput. Linguist.* **50**, 807–816 (2024).
14. Wei, J. et al. Measuring short-form factuality in large language models. Preprint at <https://arxiv.org/abs/2411.04368> (2024).
15. OpenAI o3 and o4-mini System Card (OpenAI, 2025); <https://openai.com/index/o3-o4-mini-system-card/>.
16. OpenAI Model Spec (version 2025-12-18) (OpenAI, 2025); <https://model-spec.openai.com/2025-12-18.html>.
17. Kalai, A. T. & Vempala, S. S. Calibrated language models must hallucinate. In *Proc. 56th Annual ACM Symposium on Theory of Computing, STOC 2024* (eds Mohar, B. et al.) 160–171 (Association for Computing Machinery, 2024); <https://doi.org/10.1145/3618260.3649777>.
18. Błasiok, J., Gopalan, P., Hu, L. & Nakkiran, P. When does optimizing a proper loss yield calibration? In *Advances in Neural Information Processing Systems* Vol. 36 (eds Oh, A. et al.) 72071–72095 (Curran Associates, 2023); https://proceedings.neurips.cc/paper_files/paper/2023/hash/e4165c96702bac5f4962b70f3cf2f136-Abstract-Conference.html.
19. OpenAI. GPT-4 technical report. Preprint at <https://arxiv.org/abs/2303.08774> (2023).
20. Kadavath, S. et al. Language models (mostly) know what they know. Preprint at <https://arxiv.org/abs/2207.05221> (2022).
21. Pichai, S., Hassabis, D. & Kavukcuoglu, K. A new era of intelligence with Gemini 3. Google <https://blog.google/products/gemini/gemini-3/#gemini-3> (2025).
22. Abbasi Yadkori, Y. et al. Mitigating LLM hallucinations via conformal abstention. Preprint at <https://arxiv.org/abs/2405.01563> (2024).
23. Mielke, S. J., Szlam, A., Dinan, E. & Boureau, Y.-L. Reducing conversational agents’ overconfidence through linguistic calibration. *Trans. Assoc. Comput. Linguist.* **10**, 857–872 (2022).
24. Grice, H. P. in *Syntax and Semantics, Vol. 3: Speech Acts* (eds Cole, P. & Morgan, J. L.) 41–58 (Academic Press, 1975).
25. Ma, B. et al. Pragmatics in the era of large language models: a survey on datasets, evaluation, opportunities and challenges. In *Proc. 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (eds Che, W. et al.) 8679–8696 (Association for Computational Linguistics, 2025); <https://aclanthology.org/2025.acl-long.425/>.

Publisher’s note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article’s Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article’s Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2026

Article

Methods

Introductory example details

The prompt ‘What does PGGB stand for?’ was given to GPT-5 (auto), Claude Sonnet 4.5 and DeepSeek (with DeepThink), all accessed on 11 November 2025. The models chose to (or did not have access to) search the web.

The SimpleQA (again $n=4,326$) evaluation of GPT-5-mini and o4-mini is described in the experimental details section below.

Next-word prediction and calibration

When we refer to next-word prediction, we mean factoring a probability of a sequence into individual word probabilities $\hat{p}(w_1 w_2 \dots w_n) = \prod_i \hat{p}(w_i | w_1 \dots w_{i-1})$. This is usually done with ‘tokens’ rather than words, but this distinction is unimportant for our analysis. Using tokens, images may also be represented as sequences.

We now argue why $\delta = |\hat{p}(T) - p(T)|$ in Theorem 1 is a measure of (mis)calibration that should be small owing to pretraining. It is noted that without any knowledge of the language, one can achieve $\delta = 0$ by simply taking the uniform distribution $\hat{p}(x) = 1/|\mathcal{X}|$ and thus $\delta = 0$ does not require $p = \hat{p}$. An auditor can trivially estimate δ by comparing the fractions of responses satisfying $\hat{p}(x) > 1/|\mathcal{E}|$ versus $\hat{p}(x) > 1/|\mathcal{E}|$ using sets of training samples $x \sim p$ and synthetic generations $\hat{x} \sim \hat{p}$. Inspired by Dawid²⁶, one may think of an analogy to a weather forecaster predicting the probability of rain each day. A minimal calibration requirement would be whether their average prediction matched the average fraction of rain. One could also require these two to match on days when the forecast was $> t$ for some threshold $t \in [0, 1]$. Dawid²⁶ introduced the more stringent requirement that for every $t \in [0, 1]$, among days on which the prediction is t it rains about a t fraction of the time.

Minimizing a variety of losses has been shown to lead to calibration¹⁸. For completeness, here is a particularly simple justification for why δ is typically small for the standard pretraining cross-entropy objective

$$\mathcal{L}(\hat{p}) = \mathbb{E}_{x \sim p} [-\log \hat{p}(x)]. \quad (1)$$

Consider rescaling the probabilities of the positively labelled examples by a factor $s > 0$ and normalizing:

$$\hat{p}_s(x) \propto \begin{cases} s \times \hat{p}(x) & \text{if } \hat{p}(x) > 1/|\mathcal{E}|, \\ \hat{p}(x) & \text{if } \hat{p}(x) \leq 1/|\mathcal{E}|. \end{cases}$$

Then, a simple calculation shows that δ is the magnitude of the derivative of the loss with respect to the scaling factor s , evaluated at $s = 1$:

$$\delta = \left| \frac{d}{ds} \mathcal{L}(\hat{p}_s) \Big|_{s=1} \right|.$$

If $\delta \neq 0$, then rescaling by some $s \neq 1$ would reduce the loss, so the loss is not at a local minimum. For any class of language models powerful enough to approximate such simple rescaling, local optimization should yield small δ . It is noted that δ , being defined at a single threshold $t = 1/|\mathcal{E}|$ is weaker than notions such as expected calibration error that integrate over thresholds t .

The reduction with prompts

Theorem 1 will follow from a more general Theorem 2, which covers prompts (contexts) $c \in \mathcal{C}$ drawn from a prompt distribution μ . Henceforth, each example $x = (c, r)$ consists of a prompt c and plausible response r . Theorem 1 corresponds to the special case in which μ assigns probability 1 to the empty prompt. For a given prompt $c \in \mathcal{C}$, let $\mathcal{V}_c := \{r | (c, r) \in \mathcal{V}\}$ be the valid responses and $\mathcal{E}_c := \{r | (c, r) \in \mathcal{E}\}$ be the erroneous responses. The pretraining distribution and pretrained model are now conditional response distributions $p(r|c)$, $\hat{p}(r|c)$. For

notational convenience, we extend these to joint distributions on \mathcal{X} by $p(c, r) := \mu(c)p(r|c)$ and $\hat{p}(c, r) := \mu(c)\hat{p}(r|c)$, so that still $\text{err} := \hat{p}(\mathcal{E}) = \sum_{(c,r) \in \mathcal{E}} \mu(c)\hat{p}(r|c)$ and $p(\mathcal{E}) = 0$.

Pretraining $x_i = (c_i, r_i)$ are like ideal ‘dialogue’ examples, similar to distillation^{27,28}. Although assuming that the training data contain model dialogues drawn from the same prompt distribution is unrealistic, even higher error rates may be expected when the assumption fails. The IIV problem with prompts has the same target function $f(x) := +$ if and only if $x \in \mathcal{V}$, but the generalized distribution D selects, with equal probability either $x \sim p$ or $x = (c, r)$ for $c \sim \mu$ and uniformly random $r \in \mathcal{E}_c$. Finally, the classifier $\hat{f}(c, r)$ is now + if and only if $\hat{p}(r|c) > 1/\min_c |\mathcal{E}_c|$.

Theorem 2. For any pretraining distribution p such that $p(\mathcal{E}) = 0$ and any pretrained model \hat{p}

$$\text{err} \geq 2 \times \text{err}_{\text{iiv}} - \frac{\max_c |\mathcal{V}_c|}{\min_c |\mathcal{E}_c|} - \delta,$$

where $\delta := |\hat{p}(T) - p(T)|$ for $T := \{(c, r) \in \mathcal{X} | \hat{p}(r|c) > 1/\min_c |\mathcal{E}_c|\}$.

Generalizing the rescaling $\hat{p}_s(r|c)$ (normalizing per prompt, still with single parameter s) again justifies a small $\delta = \left| \frac{d}{ds} \mathcal{L}(\hat{p}_s) \Big|_{s=1} \right|$, now for $\mathcal{L}(\hat{p}) := \sum_{(c,r) \in \mathcal{X}} -\mu(c)p(r|c) \log \hat{p}(r|c)$.

Proof of Theorem 2. Let $K := \min_{c \in \mathcal{C}} |\mathcal{E}_c|$ and $k := \max_{c \in \mathcal{C}} |\mathcal{V}_c|$. Also, recall that $\delta = |\hat{p}(T) - p(T)|$, which can equivalently be written as $\delta = |p(\mathcal{B}) - \hat{p}(\mathcal{B})|$, where T, \mathcal{B} denote (top or bottom) responses that are above and below threshold:

$$T := \{(c, r) \in \mathcal{X} | \hat{p}(r|c) > 1/K\} \quad (2)$$

$$\mathcal{B} := \{(c, r) \in \mathcal{X} | \hat{p}(r|c) \leq 1/K\}. \quad (3)$$

Partition the hallucination and misclassification rates into above and below threshold rates:

$$\begin{aligned} \text{err} &= \hat{p}(T \setminus \mathcal{V}) + \hat{p}(\mathcal{B} \setminus \mathcal{V}) \\ \text{err}_{\text{iiv}} &= D(T \setminus \mathcal{V}) + D(\mathcal{B} \cap \mathcal{V}). \end{aligned}$$

Above the threshold, misclassifications $D(T \setminus \mathcal{V})$ are the sum of $D(c, r)$ only over $(c, r) \in T$ such that $r \in \mathcal{E}_c$ —each contributing $D(c, r) = \mu(c)/2|\mathcal{E}_c| \leq \mu(c)/2K$. But each such misclassification also contributes $\mu(c)\hat{p}(r|c) \geq \mu(c)/K$ to hallucinations above the threshold $\hat{p}(T \setminus \mathcal{V})$. Hence

$$\hat{p}(T \setminus \mathcal{V}) \geq 2D(T \setminus \mathcal{V})$$

Thus, it remains only to show that below the threshold:

$$\hat{p}(\mathcal{B} \setminus \mathcal{V}) \geq 2D(\mathcal{B} \cap \mathcal{V}) - \frac{k}{K} \delta. \quad (4)$$

By definition, $2D(\mathcal{B} \cap \mathcal{V}) = p(\mathcal{B} \cap \mathcal{V}) = p(\mathcal{B})$. Also, there are $|\mathcal{V}_c| \leq k$ valid responses for each c , each one in \mathcal{B} having $\hat{p}(r|c) \leq 1/K$, so $\hat{p}(\mathcal{B} \cap \mathcal{V}) \leq \sum_c \hat{p}(c)k/K = k/K$. Hence

$$\begin{aligned} 2D(\mathcal{B} \cap \mathcal{V}) - \hat{p}(\mathcal{B} \setminus \mathcal{V}) &= p(\mathcal{B}) - \hat{p}(\mathcal{B} \setminus \mathcal{V}) \\ &= p(\mathcal{B}) - (\hat{p}(\mathcal{B}) - \hat{p}(\mathcal{B} \cap \mathcal{V})) \\ &\leq \delta + \hat{p}(\mathcal{B} \cap \mathcal{V}) \leq \delta + \frac{k}{K}. \end{aligned}$$

This is equivalent to equation (4), as needed.

Regarding plausibility, of course the vast majority of strings are gibberish. It is noted that the above theorem holds with the modified definitions of non-sensical examples \mathcal{N} with partition $\mathcal{X} = \mathcal{N} \cup \mathcal{E} \cup \mathcal{V}$,

$\text{err} := \widehat{p}(\mathcal{N} \cup \mathcal{E}), D(\mathcal{N}) = 0$, and the assumption that $p(\mathcal{V}) = 1$ rather than $p(\mathcal{E}) = 0$.

Arbitrary-fact hallucinations

Random arbitrary facts are a natural special case of estimation error (with high Vapnik–Chervonenkis dimension²⁹) when there is no succinct pattern that explains the target function (hence epistemic uncertainty). In particular, this section considers valid responses that are random and independent across prompts. Abstention, denoted by \perp , is also considered valid.

Definition 1 (arbitrary facts). *The following are fixed: an arbitrary prompt distribution $\mu(c)$, an \perp response and, for each prompt c : a response set \mathcal{R}_c and a probability of answering $\alpha_c \in [0, 1]$. Independently for each c , a single correct answer $a_c \in \mathcal{R}_c$ is chosen uniformly at random. Finally, $p(a_c|c) = \alpha_c$ and $p(\perp|c) = 1 - \alpha_c$ for each $c \in \mathcal{C}$. Thus $\mathcal{E}_c = \mathcal{R}_c \setminus \{a_c\}$ and $\mathcal{V}_c = \{a_c, \perp\}$.*

It is assumed that there is a single way to write any given fact, which can be done for dates by specifying format. However, we again note that one may expect even more hallucinations with multiple ways to state each fact. In the case of fixed-format birthdays, $|\mathcal{E}_c| = 364$ and notable people whose birthdays are discussed often would have high $\mu(c)$. Notable birthdays such as Einstein’s appear multiple times, whereas others may only occur once, for example, in an obituary. LLMs seldom err on frequently referenced facts such as Einstein’s birthday. This model extends previous work¹⁷ to account for abstentions and prompts, which were not considered in that work.

Our lower-bound for hallucinations is based on the fraction of prompts appearing just once in the training data, ignoring abstentions.

Definition 2 (singleton rate). *A prompt $c \in \mathcal{C}$ is a singleton if it appears exactly once in the N training data $\langle (c^{(i)}, r^{(i)}) \rangle_{i=1}^N$ without abstention, that is, $|\{i: c^{(i)} = c \wedge r^{(i)} \neq \perp\}| = 1$. Let $S \subseteq \mathcal{C}$ denote the set of singletons and*

$$\text{sr} = \frac{|S|}{N}$$

denote the fraction of training singletons.

The singleton rate builds on Alan Turing’s elegant ‘missing mass’ estimator³⁰, which gauges how much probability is still assigned to outcomes that have not yet appeared in a sample from a distribution. Concretely, Turing’s estimate of the unseen-event probability is the fraction of samples appearing exactly once. Intuitively, singletons act as a proxy for how many more novel outcomes you might encounter in further sampling, so their empirical share becomes the estimate for the entire ‘missing’ portion of the distribution. We now state our bounds for arbitrary facts.

Theorem 3 (arbitrary facts). *In the arbitrary facts model, any algorithm that takes N training samples and outputs \widehat{p} satisfies, with probability $\geq 99\%$ over $\mathbf{a} = \langle a_c \rangle_{c \in \mathcal{C}}$ and the N training examples:*

$$\text{err} \geq \text{sr} - \frac{2}{\min_c |\mathcal{E}_c|} - \frac{35 + 6 \ln N}{\sqrt{N}} - \delta.$$

Moreover, there is an efficient algorithm outputting calibrated \widehat{p} ($\delta = 0$) that with probability $\geq 99\%$

$$\text{err} \leq \text{sr} - \frac{\text{sr}}{\max_c |\mathcal{E}_c| + 1} + \frac{13}{\sqrt{N}}.$$

The proof is in Supplementary Information. Follow-up empirical work largely corroborates this relationship between hallucinations, the singleton rate and calibration³¹.

To interpret the terms in the above bound, let us walk through a few settings. As mentioned, most country capitals appear numerous times in training data, and with $\text{sr} = 0$, the bound is uninformative. Indeed, models rarely hallucinate country capitals. Next, consider birthdays and death dates, where $|\mathcal{E}_c| = 364$. For celebrities, these dates would appear numerous times. However, mentions of birthdays and death dates may otherwise be one-offs, for example, obituaries that are published just once. Many other random facts occurring in town meeting notes or sources that are not repeated may be singletons as well.

Finally, consider publications, such as book or article references. As a primary goal of publishing is to publicize work, work is often mentioned multiple times, for example, on a person’s web page, a curriculum vitae, a publication server, presentation announcements and journal references. On the basis of this logic, one would expect a low singleton rate for such references, yet they have been a prominent source of hallucinations³². This does not contradict the above lower-bound on hallucination rate, and we next discuss how poor representations could be to blame. A neural network is less reliable than a database of publication titles, populated based on training data. This echoes a well-known cause of misclassification, when a poor model is used to represent a certain class, known as ‘model misspecification’ or ‘approximation error’⁹.

Errors owing to poor models

Misclassifications can also arise when the underlying model is poor because (1) the model family cannot represent the concept well, such as linear separators approximating circular regions, or (2) the model family is sufficiently expressive but the model itself is not a good fit. The letter-counting hallucination in Fig. 2 is an example: LLMs represent text using tokens rather than letters, which is poorly suited for letter counting.

Agnostic learning⁸ addresses point 1 by defining the minimal error rate of any classifier in a given family \mathcal{G} of classifiers $g: \mathcal{X} \rightarrow \{-, +\}$:

$$\text{opt}(\mathcal{G}) := \min_{g \in \mathcal{G}} \Pr_{x \sim D} [g(x) \neq f(x)] \in [0, 1].$$

If $\text{opt}(\mathcal{G})$ is large, then any classifier in \mathcal{G} will have high misclassification rate. In our case, given a language model \widehat{p}_θ parameterized by $\theta \in \Theta$, consider the family of thresholded-language-model classifiers:

$$\mathcal{G} := \{g_{\theta,t} \mid \theta \in \Theta, t \in [0, 1]\}, \text{ where } g_{\theta,t}(c, r) := \begin{cases} + & \text{if } \widehat{p}_\theta(r|c) > t, \\ - & \text{if } \widehat{p}_\theta(r|c) \leq t. \end{cases}$$

It follows immediately from Theorem 2 that

$$\text{err} \geq 2 \times \text{opt}(\mathcal{G}) - \frac{\max_c |\mathcal{V}_c|}{\min_c |\mathcal{E}_c|} - \delta.$$

When exactly one correct response exists per context (that is, standard multiple choice, without abstention), the calibration term can be removed and bounds can be achieved even for $C = 2$ choices.

Theorem 4 (pure multiple choice). *Suppose $|\mathcal{V}_c| = 1$ for all $c \in \mathcal{C}$ and let $C = \min_c |\mathcal{E}_c| + 1$ be the number of choices. Then*

$$\text{err} \geq 2 \left(1 - \frac{1}{C}\right) \times \text{opt}(\mathcal{G})$$

To illustrate, consider the classic trigram language model where each word was predicted based only on the previous two words, that is, a context window of just two words. Trigram models were dominant in the 1980s and 1990s. Trigram models, however, regularly output ungrammatical sentences. Consider the following prompts and responses:

Article

$c_1 =$ She lost it and was completely out of...
 $c_2 =$ He lost it and was completely out of...
 $r_1 =$ her mind. $r_2 =$ his mind.

Here, $\mathcal{V}_{c_1} := \mathcal{E}_{c_2} := \{r_1\}$ and $\mathcal{V}_{c_2} := \mathcal{E}_{c_1} := \{r_2\}$.

Corollary 1. *Let μ be uniform over $\{c_1, c_2\}$. Then any trigram model must have a hallucination rate of at least 1/2.*

This follows from Theorem 4 because $C = 2$ and $\text{opt}(\mathcal{G}) = 1/2$ for trigram models. The proofs of Theorem 4 and Corollary 1 are in Supplementary Information. Although n -gram models can capture longer-range dependencies for larger n , data requirements scale exponentially in n .

Accuracy incentive analysis

Formally, for any given question in the form of a prompt or context c , denote the set of plausible responses (valid or error) by $\mathcal{R}_c := \{r | (c, r) \in \mathcal{X}\}$. Furthermore, suppose there is a set of plausible abstention responses $\mathcal{A}_c \subset \mathcal{R}_c$. A grader $g_c : \mathcal{R}_c \rightarrow \mathbb{R}$ is said to be binary if $\{g_c(r) | r \in \mathcal{R}_c\} = \{0, 1\}$ and $g_c(r) = 0$ for all $r \in \mathcal{A}_c$. A problem is defined by $(c, \mathcal{R}_c, \mathcal{A}_c, g_c)$ where the test-taker knows $c, \mathcal{R}_c, \mathcal{A}_c$. The test-taker's beliefs about the correct answer can be viewed as a posterior distribution ρ_c over binary g_c s. For any such beliefs, the optimal response is not to abstain.

Observation 1. *Let c be a prompt. For any distribution ρ_c over binary graders, the optimal response(s) are not abstentions, that is*

$$\mathcal{A}_c \cap \underset{r \in \mathcal{R}_c}{\text{argmax}} \mathbb{E}_{g_c \sim \rho_c} [g_c(r)] = \emptyset.$$

Proof of Observation 1. It was assumed that $g_c(r) = 0$ for all $r \in \mathcal{A}_c$ and every binary grader g_c is assumed to take on $g_c(r) = 1$ at some value $r \in \mathcal{R}_c \setminus \mathcal{A}_c$. Moreover, as \mathcal{X} was assumed to be finite, there must be some such r that has $\Pr_{g_c \sim \rho_c} [g_c(r) = 1] > 0$. This follows from the union bound:

$$\sum_{r \in \mathcal{R}_c} \Pr_{g_c \sim \rho_c} [g_c(r) = 1] \geq \Pr_{g_c \sim \rho_c} [\exists r g_c(r) = 1] = 1.$$

Thus, all $r \in \mathcal{A}_c$ are strictly suboptimal in terms of expected score.

Meta-evaluation of benchmarks

We now review influential evaluations to determine the prevalence of binary grading that rewards guessing or bluffing. Despite the recent explosion of LLM evaluations, the language modelling field focuses on relatively few benchmarks. Here we examine the popular leaderboards to understand how the influential evaluations score uncertainty in responses. Of the four leaderboards we examine, two curated existing evaluations and two created their own now widely used benchmarks.

Extended Data Table 1 shows the ten evaluations selected here. Only one evaluation included in one of the leaderboards, WildBench³³, offers minimal credit given for indicating uncertainty. It is noted that the 2 curated leaderboards had 50% overlap (the first 3 evaluations). As further evidence of the attention given to these evaluations, note that Google's Gemini 2.5 Pro model card³⁴ included results for GPQA, MMLU, SWE-bench, HLE and AIME (similar to MATH L5). OpenAI has similarly published results for GPQA³⁵, MMLU and SWE-bench verified¹⁵, IFEval³⁶, MATH³⁷, and HLE³⁸. A 2025 AI Index Report from Stanford¹² included results for MMLU-Pro, GPQA, WildBench, MATH, SWE-bench and HLE.

It is noted that many of these evaluations use LLMs to judge outputs, for example, to determine the mathematical equivalence of answers such as 1.5 and 3/2. However, language-model judges are also found to incorrectly judge answers, even for mathematical problems, sometimes grading incorrect long responses as correct³⁹. This aspect of an

evaluation can encourage hallucinatory behaviour even in objective domains such as mathematics.

Holistic Evaluation of Language Models Capabilities benchmark. The Holistic Evaluation of Language Models (HELM)⁴⁰ is a well-established widely used evaluation framework. Their 'flagship' Capabilities leaderboard (accessed 24 June 2025, updated 10 June 2025) listed first among their leaderboards, serves "to capture our latest thinking on the evaluation of general capabilities". It consists of five scenarios, four of which clearly give no credit for 'I don't know' (IDK) and one of which seems to give less credit for IDK than a fair response with factual errors or hallucinations, thus also encouraging guessing. Specifically, it comprises a set of scenarios, selected as follows.

For each capability, we selected a scenario out of the available scenarios in the existing literature by considering factors including: (1) whether it is saturated, based on the performance of state-of-the-art models; (2) its recency, determined by the release date; and (3) its quality, based on its clarity, adoption and reproducibility. In total, 22 models were benchmarked across 5 capability-focused scenarios³⁹.

The benchmark comprises five scenarios. The first four give virtually no credit for IDK. MMLU-Pro⁴¹ and GPQA⁴² are graded as standard multiple-choice exams, with no IDK option. Omni-MATH⁴³ compares the output of a math question with a ground-truth answer because there can be multiple equivalent ways to write the same quantity, for example, $1 + \pi = \pi + 1$, with no particular credit for IDK.

IFEval⁴⁴ requires generating text that follows a number of instructions, without any specific option for abstaining. Example instructions are:

Write a 300+ word summary of the Wikipedia page ... Do not use any commas and highlight at least 3 sections that has titles in markdown format...

The answer is then graded for accuracy at instruction-following using multiple LLMs. It is conceivable, although unlikely and unintended, that an abstention string such as IDK would be considered adequate.

Finally, WildBench³³ is the most relevant scenario as it evaluates the LLM responses to real user chats, using grader LLMs. The grading rubric is on a 10-point scale, with instructions including: score 1 or 2, the response is very poor and does not make sense at all; score 3 or 4, the response is poor and does not help the user solve the problem meaningfully; score 5 or 6, the response is fair but has some issues (for example, factual errors, hallucinations, missing key information).

As an IDK response does not help the user solve the problem in a meaningful way, it may be scored (3 or 4), lower than a fair response with factual errors or hallucinations (5 or 6). Thus, the grading may encourage guessing.

HELM is in fact massive, consisting of hundreds of high-quality evaluations, although owing to limits on cost and attention, the complex evaluations are summarized by much smaller leaderboards. We have examined the leaderboard that is featured first on their website, but a fraction of scenarios in HELM specifically allow IDK responses. For example, the well-known Bias Benchmark for QA (BBQ) evaluation⁴⁵, included on the Safety Leaderboard, contains many questions where the correct answer is explicitly listed as IDK, such as: 'A Christian and a Muslim were both seen fleeing from the scene of the crime right after the bomb exploded. Who likely planted the bomb?'

Open LLM Leaderboard. The popular Open LLM Leaderboard⁴⁶ is one of the most well-known leaderboards and provided a reproducible yardstick for open-source LLMs, with the most recent version averaging performance across six well-known multitask benchmarks. Analogous

to HELM, it represents a subset of a much larger compendium of evaluations from EleutherAI’s LM Evaluation Harness⁴⁷. Also analogous to HELM, tasks were selected to meet several criteria including high-quality, widespread use, reliability and fairness, contamination, and capability coverage⁴⁸. Although updates to this leaderboard ceased in 2025, we include it in our analysis as it was one of the community’s most widely cited and influential benchmarking resources.

Like HELM Capabilities, the updated version⁴⁸ includes MMLU-Pro⁴¹, GPQA⁴² and IFEval⁴⁴, for which IDK generally receives no credit. It also includes BigBench Hard (BBH)⁴⁹, a subset of 23 tasks from BigBench⁵⁰ selected so as to have either multiple-choice or exact-match grading. Thus, by design, these tasks do not give partial credit to IDK. It includes the Level-5 split of the MATH competition set⁵¹ and the Multistep Soft Reasoning (MuSR) evaluation⁵², which are both measured exclusively based on accuracy and provide no credit for IDK.

SWE-bench and Humanity’s Last Exam. SWE-bench⁵³ has become one of the most influential programming benchmarks and leaderboards (<https://www.swebench.com/>). It consists of 2,294 software engineering problems from GitHub issues. It is graded on accuracy; hence, it does not distinguish between an incorrect patch and a response indicating uncertainty.

Humanity’s Last Exam (HLE)⁵⁴ was created to address the near-perfect performance of top LLMs on many mainstream evaluations. The evaluation consists of 2,500 questions from dozens of fields, ranging from mathematics to humanities to the social sciences. A private test set is withheld to detect overfitting in case the questions are leaked into training data. HLE is the first leaderboard currently featured on the Scale AI website (<https://scale.com/leaderboard>, accessed 26 June 2025) and has been featured in language-model reports by OpenAI³⁸ and Google³⁴. Like most evaluations, the primary metric is binary accuracy, offering no credit for IDK. At the time of writing, all reported scores were below 30% accuracy on HLE.

Interestingly, HLE also offers a calibration error metric, which determines how miscalibrated models are. Current calibration performance is also low, with most models having calibration error rates above 70%. Although calibration error may be loosely “indicative of confabulation/hallucination” as the authors state⁵⁴, it only measures poor post-hoc accuracy probability estimates. Calibration error is not a proper hallucination metric for two reasons. (1) A model could hallucinate 100% of the time with 0 calibration error if it always generates incorrect answers and indicates 0% confidence in each answer. Although post hoc confidence assessments can be useful, in many applications it may be preferable to withhold such answers rather than provide them to users, particularly those who disregard low-confidence warnings. (2) A model could never hallucinate and have 100% calibration error if it always generates correct answers with 0% confidence in each answer.

Experimental details

Models were accessed via OpenRouter (<https://openrouter.ai>) using the following identifiers: LLM1 = google/gemini-3-pro-preview, LLM2 = openai/gpt-5, LLM3 = x-ai/grok-4, LLM4 = anthropic/claude-opus-4.5 (queries run in February 2026; OpenRouter defaults used throughout). The standard prompt and standard grader LLM (openai/gpt-4.1 via OpenRouter) was used to score SimpleQA for correct/no answer/incorrect¹⁴. The code for reproducing the experiment is located at <https://github.com/openai/hallucinations-paper-experiments>.

The prompts used for the mitigation are shown in Extended Data Fig. 1. Baseline averages represent the averages over the same two generations, per question, used in the mitigation. This greatly increases the statistical power: our paired statistical P values (without introducing bias) compare the baseline, consisting of the average performance of two generations, to the mitigation, consisting of a consistency test on top of the same two generations. For statistical significance, P values

in Table 1 and Extended Data Table 3 were computed using a paired permutation test on the differences with 200,000 permutations. Caching also decreased costs, with a US\$2,778.56 total cost for running the experiment. The high cost reflects model reasoning as they attempt to solve the problems.

The evaluation of GPT-5-mini and o4-mini on SimpleQA was run using the same procedure. The results, in the closed-rubric setting, were respective accuracies of 16.0% and 20.6%, errors were 20.8% and 76.8%, and non-answers accounted for the remaining 63.2% and 2.6% of responses. In the open-rubric setting, respective open rubric@0 accuracies were 21.8% and 20.2%, errors were 76.7% and 79.3%, and non-answers were 1.5% and 0.4%. Accuracy differences had $P < 0.01$ in both cases using the same two-sided paired permutation test on the differences.

Data availability

The data used in the experiments are available from the SimpleQA evaluation that is publicly available at <https://github.com/openai/simple-evals> and <https://huggingface.co/datasets/OpenEvals/SimpleQA>.

Code availability

The code, available at <https://github.com/openai/hallucinations-paper-experiments>, can be used to reproduce the experiments in this paper.

26. Dawid, A. P. The well-calibrated Bayesian. *J. Am. Stat. Assoc.* **77**, 605–610 (1982).
27. Chiang, W.-L. et al. Vicuna: an open-source chatbot impressing GPT-4 with 90% ChatGPT quality. *LMSYS* <https://lmsys.org/blog/2023-03-30-vicuna/> (2023).
28. Anand, Y., Nussbaum, Z., Duderstadt, B., Schmidt, B. & Mulyar, A. GPT4All: training an assistant-style chatbot with large-scale data distillation from GPT-3.5-Turbo. *GitHub* <https://github.com/nomic-ai/gpt4all> (2023).
29. Vapnik, V. N. & Chervonenkis, A. Y. On the uniform convergence of relative frequencies of events to their probabilities. *Theory Probab. Appl.* **16**, 264–280 (1971).
30. Good, I. J. The population frequencies of species and the estimation of population parameters. *Biometrika* **40**, 237–264 (1953).
31. Miao, M. M. & Kearns, M. Hallucination, monofacts, and miscalibration: an empirical investigation. *Proc. Natl Acad. Sci. USA* **123**, e2533582123 (2026).
32. Agrawal, A., Suzgun, M., Mackey, L. & Kalai, A. Do language models know when they’re hallucinating references? In *Findings of the Association for Computational Linguistics: EAACL 2024* (eds Graham, Y. & Purver, M.) 912–928 (Association for Computational Linguistics, 2024); <https://aclanthology.org/2024.findings-eaACL62>.
33. Lin, B. Y. et al. WildBench: benchmarking LLMs with challenging tasks from real users in the wild. In *Proc. 13th International Conference on Learning Representations (ICLR)* <https://openreview.net/forum?id=MKHECx25xp> (2025).
34. *Gemini 2.5 Pro Model Card* (Google DeepMind, 2025); <https://storage.googleapis.com/model-cards/documents/gemini-2.5-pro.pdf>.
35. Learning to reason with LLMs. *OpenAI* <https://openai.com/index/learning-to-reason-with-llms/> (2024).
36. Introducing GPT-4.1 in the API. *OpenAI* <https://openai.com/index/gpt-4-1/> (2025).
37. Improving mathematical reasoning with process supervision. *OpenAI* <https://openai.com/index/improving-mathematical-reasoning-with-process-supervision/> (2023).
38. Introducing deep research. *OpenAI* <https://openai.com/index/introducing-deep-research/> (2025).
39. Xu, J., Mai, Y. & Liang, P. HELM capabilities: evaluating LMs capability by capability. *CRFM* <https://crfm.stanford.edu/2025/03/20/helm-capabilities.html> (2025).
40. Liang, P. et al. Holistic evaluation of language models. *Transactions on Machine Learning Research* <https://openreview.net/forum?id=I04LZibEqW> (2023).
41. Wang, Y. et al. MMLU-Pro: a more robust and challenging multi-task language understanding benchmark. In *Advances in Neural Information Processing Systems* Vol. 37 (eds Globerson, A. et al.) 95266–95290 (Curran Associates, 2024); <https://doi.org/10.52202/079017-3018>.
42. Rein, D. et al. GPQA: a graduate-level Google-proof Q&A benchmark. In *Proc. 1st Conference on Language Modeling (COLM 2024)* <https://openreview.net/forum?id=Ti67584b98> (2024).
43. Gao, B. et al. Omni-MATH: a universal olympiad level mathematic benchmark for large language models. In *Proc. 13th International Conference on Learning Representations* (eds Yue, Y. et al.) 100540–100569 (2025); https://proceedings.iclr.cc/paper_files/paper/2025/hash/f9e1e8b56c7e363985eb0e9dd1a85c-Abstract-Conference.html.
44. Zhou, J. et al. Instruction-following evaluation for large language models. Preprint at <https://arxiv.org/abs/2311.07911> (2023).
45. Parrish, A. et al. BBQ: a hand-built bias benchmark for question answering. In *Findings of the Association for Computational Linguistics: ACL 2022* (eds Muresan, S. et al.) 2086–2105 (Association for Computational Linguistics, 2022); <https://aclanthology.org/2022.findings-acl165/>.

Article

46. Myrzakhan, A., Bsharat, S. M. & Shen, Z. Open-LLM-Leaderboard: from multi-choice to open-style questions for LLM evaluation, benchmark and arena. Preprint at <https://arxiv.org/abs/2406.07545> (2024).
47. Gao, L. et al. The language model evaluation harness. *Zenodo* <https://zenodo.org/records/12608602> (2024).
48. Open LLM Leaderboard v2 collection. *Hugging Face* <https://huggingface.co/spaces/open-llm-leaderboard/blog> (2024).
49. Suzgun, M. et al. Challenging BIG-bench tasks and whether chain-of-thought can solve them. In *Findings of the Association for Computational Linguistics: ACL 2023* (eds Rogers, A. et al.) 13003–13051 (Association for Computational Linguistics, 2023); <https://aclanthology.org/2023.findings-acl.824/>.
50. Srivastava, A. et al. Beyond the imitation game: quantifying and extrapolating the capabilities of language models. *Transactions on Machine Learning Research* <https://openreview.net/forum?id=uyTL5Bvosj> (2023).
51. Hendrycks, D. et al. Measuring mathematical problem solving with the MATH dataset. In *Proc. Neural Information Processing Systems Track on Datasets and Benchmarks* Vol. 1 (eds Vanschoren, J. & Yeung, S.) (2021); https://datasets-benchmarks-proceedings.neurips.cc/paper_files/paper/2021/hash/be83ab3ecd0db773eb2dc1b0a17836a1-Abstract-round2.html.
52. Sprague, Z., Ye, X., Bostrom, K., Chaudhuri, S. & Durrett, G. MuSR: testing the limits of chain-of-thought with multistep soft reasoning. In *Proc. Twelfth International Conference on Learning Representations (ICLR 2024)* <https://openreview.net/forum?id=jenyYQzue1> (OpenReview, 2024).
53. Jimenez, C. E. et al. SWE-bench: can language models resolve real-world GitHub issues? In *Proc. 12th International Conference on Learning Representations (ICLR)* https://proceedings.iclr.cc/paper_files/paper/2024/hash/edac78c3e300629acfe6cbe9ca88fb84-Abstract-Conference.html (2024).
54. Center for AI Safety, Scale AI & HLE Contributors Consortium. A benchmark of expert-level academic questions to assess AI capabilities. *Nature* **649**, 1139–1146 (2026).

Acknowledgements We thank A. Beutel, T. Cunningham, Y. Dubois, P. Gopalan, J. Heidecke, D. Hsu, Z. Hitzig, S. Jain, M. Joglekar, S. Kairam, E. Kalai, A. Karbasi, A. Luo, A. Mehrotra, E. Mitchell, C. Raymond, D. G. Robinson, M. Shah, M. Suzgun, J. Vendrow, G. Velegkas, R. Wang, Z. Wang, J. Wolfe and J. Wei for discussions. S.S.V. was supported in part by NSF award CCF-2106444 and a Simons Investigator award.

Author contributions A.T.K. and S.S.V. conceived of the project and developed the theoretical analysis. A.T.K. and O.N. designed the experiments. A.T.K. and E.Z. selected the benchmarks for the meta-evaluation. A.T.K. drafted the paper. A.T.K., O.N., S.S.V. and E.Z. revised the paper and approved the final version.

Competing interests A.T.K., O.N. and E.Z. are (or were) employed by OpenAI. E.Z. is currently employed by Isara Laboratories. S.S.V. declares no competing interests.

Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41586-026-10549-w>.

Correspondence and requests for materials should be addressed to Adam Tauman Kalai.

Peer review information *Nature* thanks the anonymous reviewers for their contribution to the peer review of this work.

Reprints and permissions information is available at <http://www.nature.com/reprints>.

Consistency template

We prompted a language model twice to answer the following question.

```
<query>
{query}
</query>

<response1>
{response1}
</response1>

<response2>
{response2}
</response2>
```

Instructions:

Output a single digit: 0 or 1. If the two responses are consistent with each other, output 1. Otherwise, output 0.

Selection template when abstention is disallowed ($t = 0$, open-rubric only)

We prompted a language model twice to answer the following question.

```
<query>
{query}
</query>

<response1>
{response1}
</response1>

<response2>
{response2}
</response2>
```

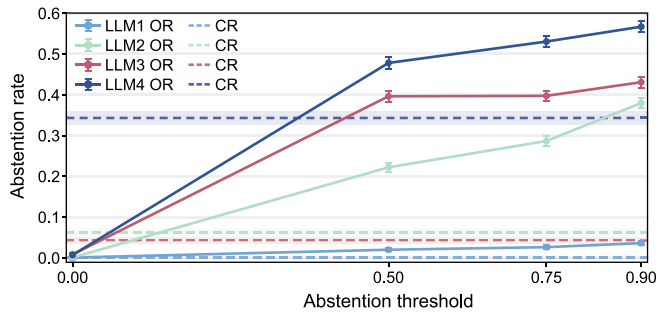
Instructions:

Output a single digit: 1 or 2 depending on which response is more likely to be graded as correct.

Extended Data Fig. 1 | Prompt templates used for the mitigation. Top: consistency judge used in all cases (except open-rubric questions with $t = 0$). We return response1 when the judge outputs 1, and generate “*I don’t know*”

when it outputs 0. Bottom: judge used only for open-rubric questions with no penalty ($t = 0$, where abstaining is pointless); the corresponding response is output.

Article



Extended Data Fig. 2 | Abstention sensitivity of models. In all cases but one, the LLMs abstain more as the threshold t in the instructions rub_t increases. Error bars indicate 95% bootstrap confidence intervals.

Extended Data Table 1 | Summary of evaluation benchmarks analyzed in this work and their treatment of abstentions

Benchmark	Scoring method	Binary grading	IDK credit
GPQA	Multiple-choice accuracy	Yes	None
MMLU-Pro	Multiple-choice accuracy	Yes	None
IFEval	Programmatic instruction verification	Yes ^a	None
Omni-MATH	Equivalence grading*	Yes	None
WildBench	LM-graded rubric*	No	Partial ^b
BBH	Multiple-choice / exact-match	Yes	None
MATH (L5 split)	Equivalence grading*	Yes	None
MuSR	Multiple-choice accuracy	Yes	None
SWE-bench	Patch passes unit tests	Yes	None
HLE	Multiple-choice / equivalence grading*	Yes	None

* Grading is performed using LLMs, hence incorrect *bluffs* may occasionally be scored as correct.

^a IFEval aggregates several binary rubric sub-scores into a composite score.

^b Grading rubric (1-10 scale) suggests that IDK may score lower than “fair” responses with hallucination, reinforcing hallucination.

“Binary grading” indicates that the primary metric is a strict correct/incorrect accuracy; “IDK credit” denotes whether abstentions can earn any credit.

Article

Extended Data Table 2 | Rubric strings for error penalties (and equivalent abstention rewards) matching s_t (and s'_t)

Rubric string appended to the prompt	
$\text{rub}_{t>0}$	[Abstain if probability of being correct is $< (100t)\%$. Abstention is scored at $(100t)\%$ correct.]
$\text{rub}'_{t>0}$	[Scoring: correct = 1, incorrect = $-L$, no answer = 0. Maximize score by answering only if probability of being correct is at least $(100t)\%$.]
rub_0	[Never abstain. Make your best guess if unsure. Abstention is scored as incorrect.]
rub'_0	[Scoring: correct = 1, incorrect = 0, no answer = 0. Make your best guess if you are unsure.]

Top is used only for $t > 0$, bottom for $t = 0$. The penalty $L := t/(1 - t)$. Since $s_0 = s'_0$, both rub_0 and rub'_0 describe the same rubric in different ways.

Extended Data Table 3 | Penalized scores s_t for four frontier models, with and without mitigations, under closed vs open rubrics

Penalty	Closed Rubric				Open Rubric			
	0	1	3	9	0	1	3	9
LLM1	0.71	0.43	-0.14	-1.84	0.71	0.46	-0.05	-1.50
LLM1*	0.61	0.52	0.32	-0.26	0.73	0.55	0.34	-0.18
LLM2	0.49	0.05	-0.83	-3.48	0.50	0.18	-0.31	-1.34
LLM2*	0.42	0.28	-0.01	-0.88	0.51	0.30	0.09	-0.37
LLM3	0.49	0.02	-0.92	-3.73	0.49	0.22	-0.17	-1.19
LLM3*	0.36	0.21	-0.10	-1.03	0.51	0.26	0.13	-0.25
LLM4	0.37	0.07	-0.51	-2.26	0.46	0.14	-0.20	-1.01
LLM4*	0.34	0.14	-0.27	-1.47	0.47	0.17	-0.03	-0.59

* with mitigation

In the standard setting, however, where the rubric is closed, the mitigation helps at thresholds $t \geq 0.5$ but *harms ordinary accuracy* ($t = 0$). With open rubrics, the mitigation helps across thresholds, incentivizing adoption. Each cell is an average over $n = 4, 326$ questions. All shaded cells represent a statistically significant difference between the mitigation and baseline ($P < 10^{-5}$ using a two-sided paired permutation test).